

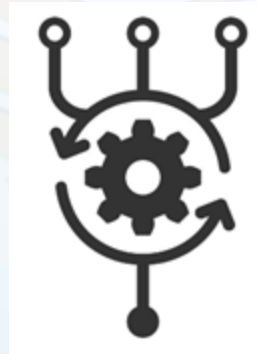
Fairness in Machine Learning - Overview

Arun Rajkumar
Dept of DSAI, IITM

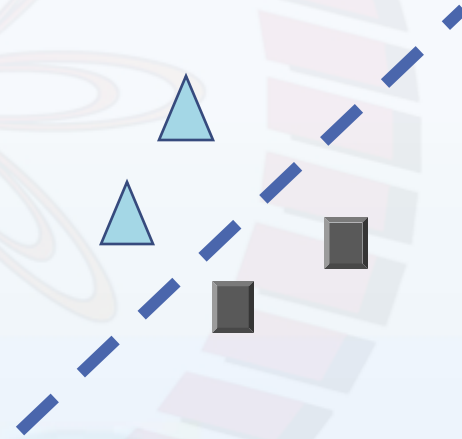
Machine Learning Pipeline



Data



Algorithm



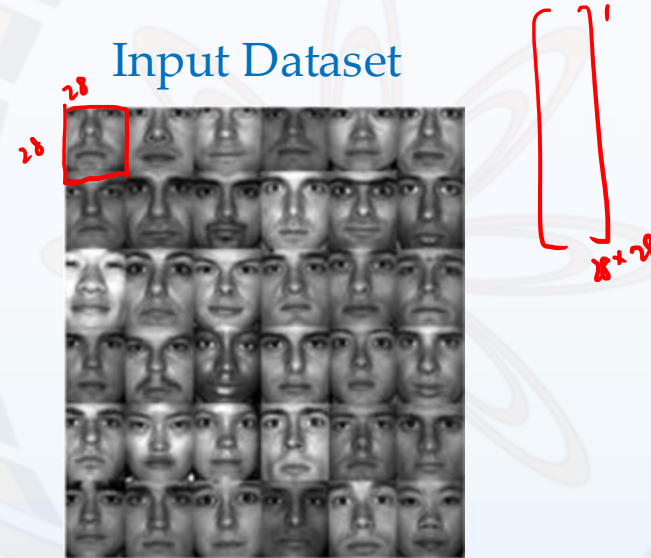
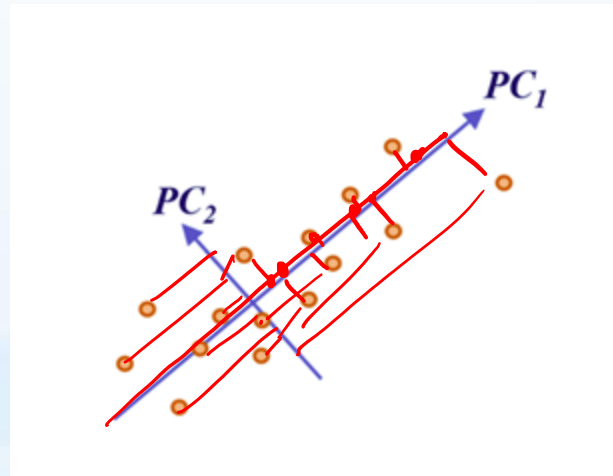
Model



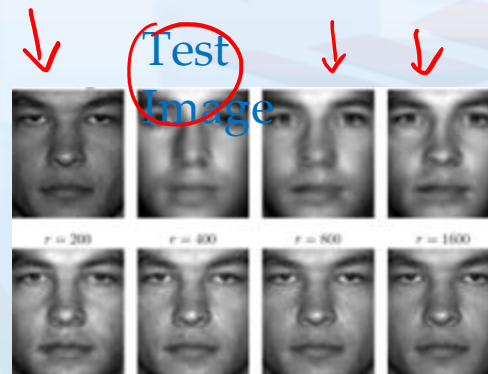


NPTTEL

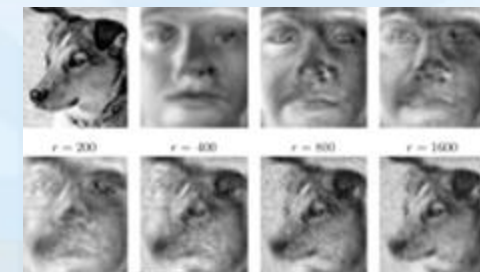
Principal Component Analysis - Recap



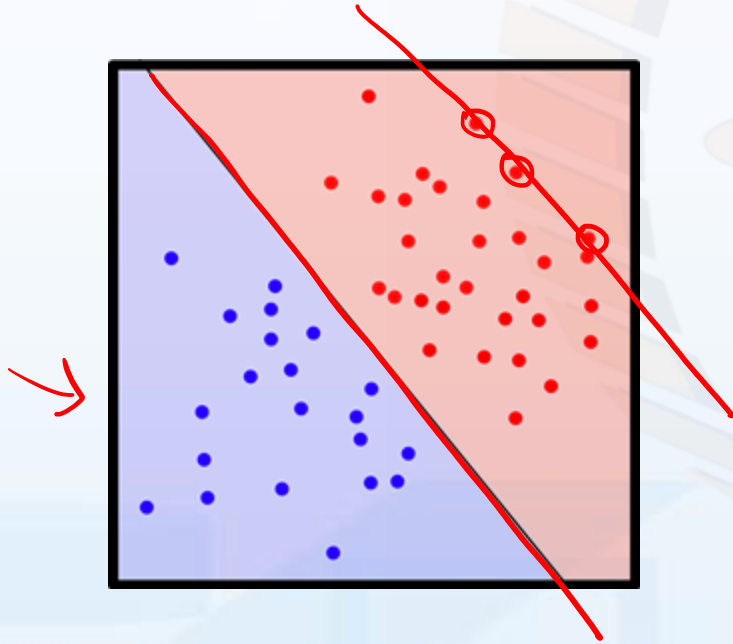
Principal Components



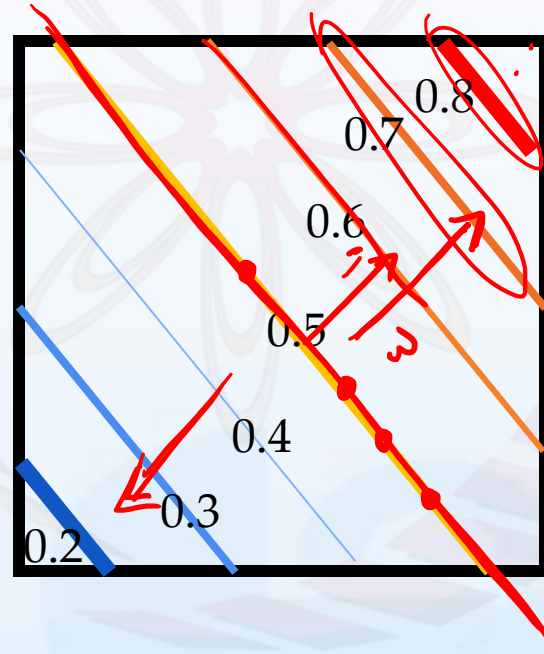
Another Test Image



Logistic Regression - Recap

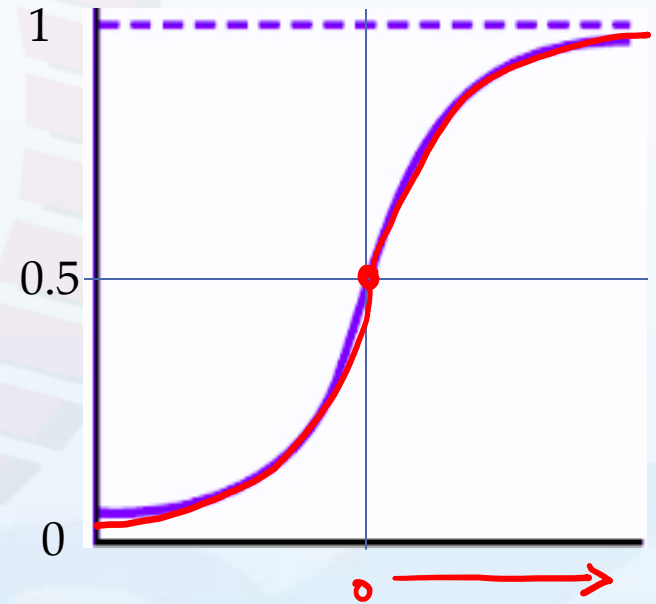


Linearly Separable Dataset



Probabilistic Model

$$P(y=+1/x)$$



Logistic Function

$$P(y=1/x) = 1/(1+\exp(-\bar{w} \cdot x))$$

Logistic Regression - Recap

$$P(y=+1/x) = g(z)$$

- $g(z) = 0.5$ if $z = 0$
- $g(z) \rightarrow 1$ as $z \rightarrow \infty$
- $g(z) \rightarrow 0$ as $z \rightarrow -\infty$

Logistic Regression - Recap

Model: LOGISTIC REGRESSION

$$P(y=1/x) = \frac{1}{1 + e^{-w^T x}}$$

Data : $\{(x_1, y_1), \dots, (x_n, y_n)\}$

- How to find w ?

- Maximum likelihood!

$$L(\underline{w}; \text{Data}) = \prod_{i=1}^n (g(w^T x_i))^{y_i} \cdot (1 - g(w^T x_i))^{(1-y_i)}$$

Logistic Regression - Recap

- No closed form solution for maximisation
- Can perform gradient descent

Gradient

$$\nabla \log L(w) = \sum_{i=1}^n x_i \left(y_i - \frac{1}{1 + e^{-w^T x_i}} \right)$$



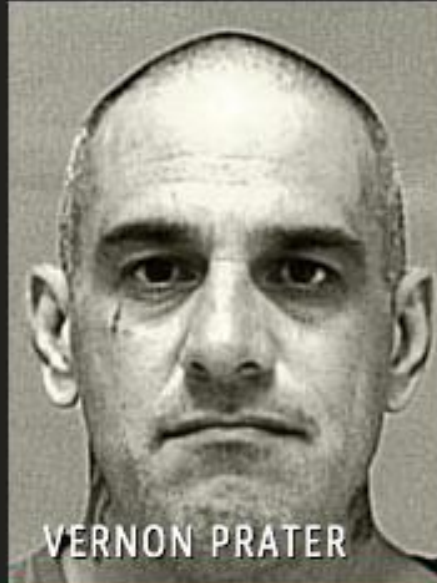
NPTTEL



Google apologises for Kannada is the ‘ugliest language in India’ search result after backlash

“Sometimes, the way content is described on the Internet can yield surprising results to specific queries. We know this is not ideal, but we take swift corrective action when we are made aware of an issue and are continually working to improve our algorithms. Naturally, these are not reflective of the opinions of Google, and we apologize for the misunderstanding and hurting any sentiments,” a Google spokesperson said in a response to Hindustan Times.

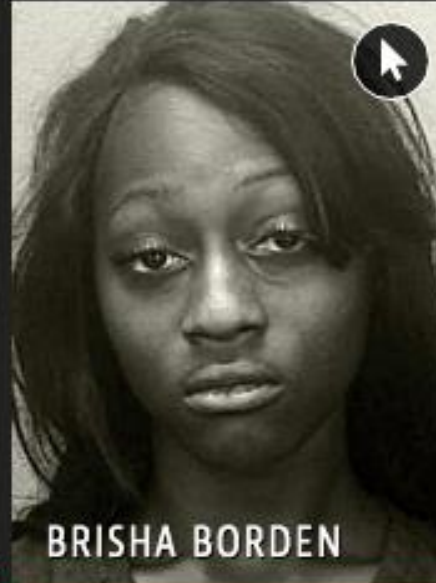
Two Petty Theft Arrests



VERNON PRATER

LOW RISK

3



BRISHA BORDEN

HIGH RISK

8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Two Petty Theft Arrests

VERNON PRATER

Prior Offenses

2 armed robberies, 1
attempted armed
robbery

Subsequent Offenses

1 grand theft

LOW RISK

3

BRISHA BORDEN

Prior Offenses

4 juvenile
misdemeanors

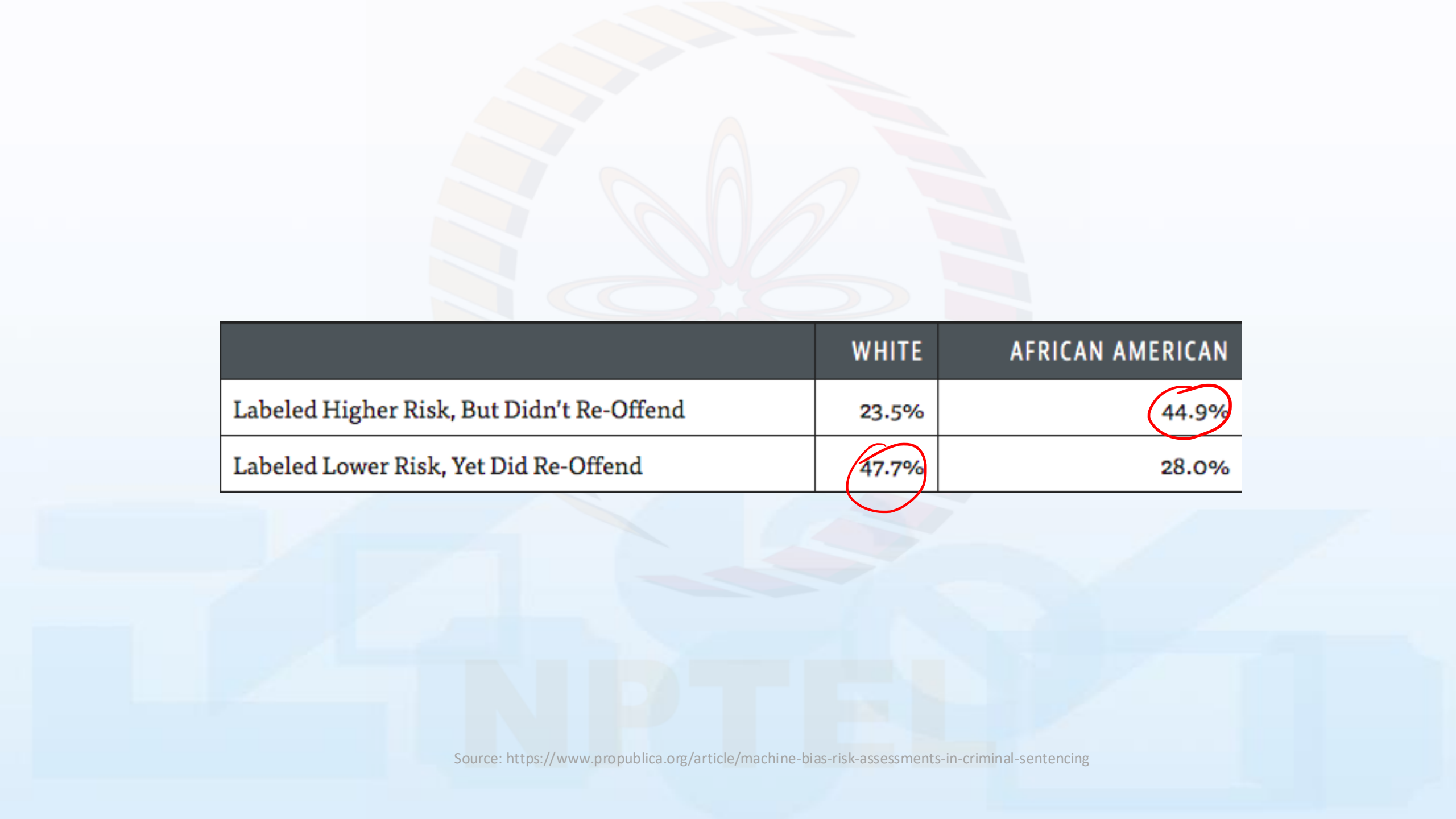
Subsequent Offenses

None

HIGH RISK

8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.



	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

Source: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>



Amazon scraps secret AI recruiting tool that showed bias against women

NPTTEL

Discrimination in Online Ad Delivery

Latanya Sweeney
Harvard University
latanya@fas.harvard.edu

January 28, 2013¹

Problem Statement

Given online searches of racially identifying names, show that associated personalized ads suggestive of an arrest record do not differ by race.

Biases in Word Embeddings

Wikipedia

Sexist prejudice

Profession		Sentiment	
Woman	Man	Woman	Man
Nurse	Officer	Wedding	Reinforcement
Secretary	Hunter	Divorce	Attack
Teacher	Commander	Anulment	Combat
Saleswoman	Guard	Engagement	Power
Actress	Cameraman	Marry	Decrease

Social Media

Sexist prejudice

Profession		Sentiment	
Woman	Man	Woman	Man
Nurse	Policeman	Agitation	Robber
Secretary	Musician	Mature	Attacker
Pharmacist	Priest	Love	Injured
Religion teacher	Coach	Increase	Fascist
Correspondent	Paramedic	Stubbornness	Overwhelmed

<https://blog.acolyer.org/2020/12/08/bias-in-word-embeddings/>

ChatGPT



Examples

"Explain quantum computing in simple terms" →

"Got any creative ideas for a 10 year old's birthday?" →

"How do I make an HTTP request in Javascript?" →



Capabilities

Remembers what user said earlier in the conversation

Allows user to provide follow-up corrections

Trained to decline inappropriate requests



Limitations

May occasionally generate incorrect information

May occasionally produce harmful instructions or biased content

Limited knowledge of world and events after 2021





How do we fix it?

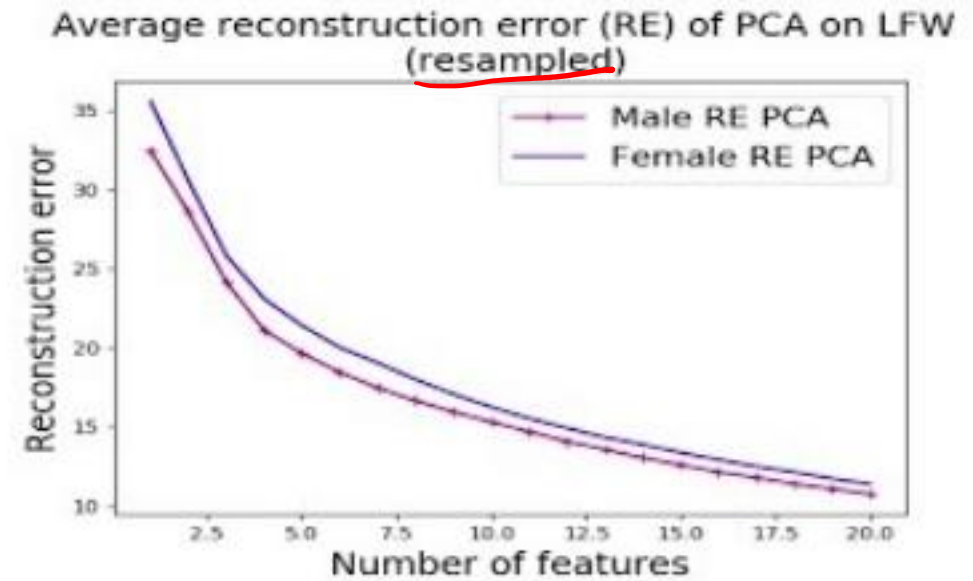
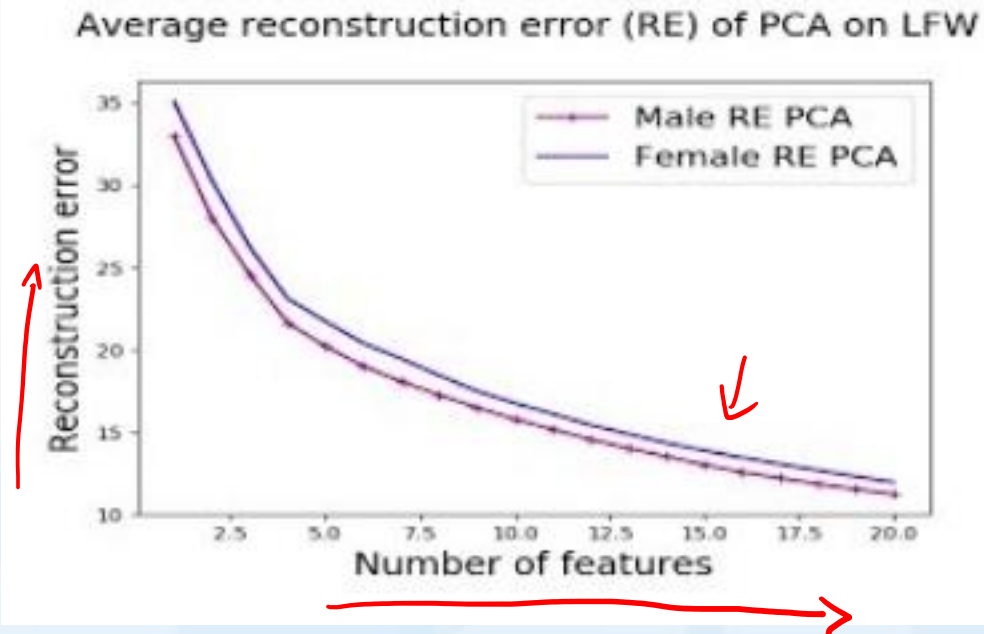
But first, what exactly is fairness?

NPTEL



Fair unsupervised learning

UnFair PCA



<https://sites.google.com/site/ssamadi/home/fair-pca-homepage>

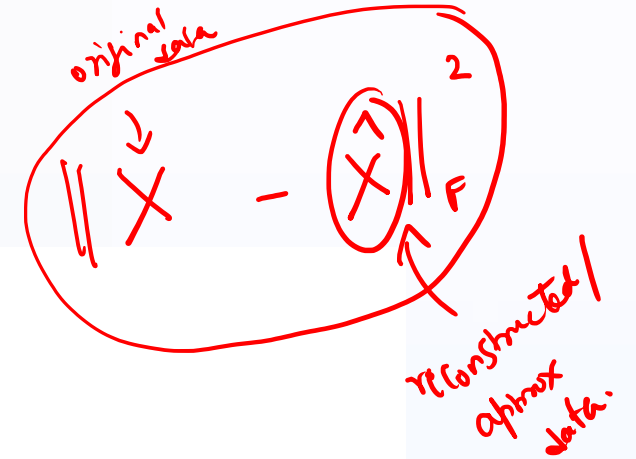
UnFair PCA

male female.

$$\begin{bmatrix} A \\ B \end{bmatrix} \rightarrow U = \begin{bmatrix} U_A \\ U_B \end{bmatrix}$$

Definition: Loss of a population A when approximated by matrix U_A . Let \hat{A} be the optimal rank- d approximation of A .

$$\text{Loss}_A(U_A) = \text{Error}_A(U_A) - \text{Error}_A(\hat{A}) = \|A - U_A\|_F^2 - \|A - \hat{A}\|_F^2$$



UnFair PCA

$$x = \begin{bmatrix} A \\ B \end{bmatrix} \rightarrow U = \begin{bmatrix} U_A \\ U_B \end{bmatrix}$$

Fair PCA problem:

$$\min_{U: \text{rank}(U)=d} \max \left\{ \frac{1}{|A|} \text{Loss}_A(U_A), \frac{1}{|B|} \text{Loss}_B(U_B) \right\}$$

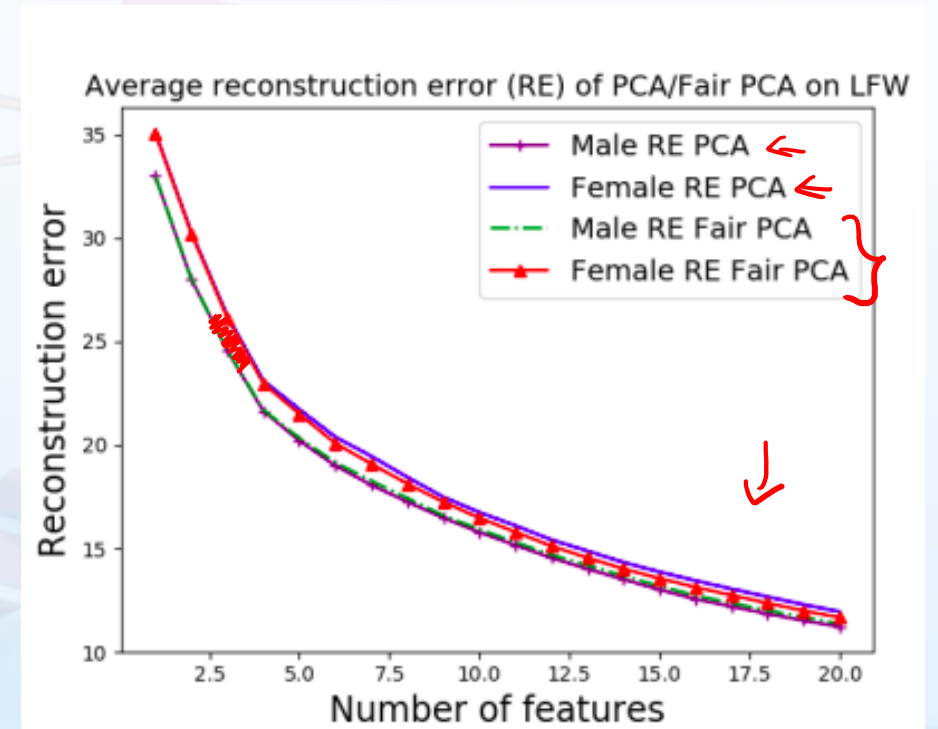
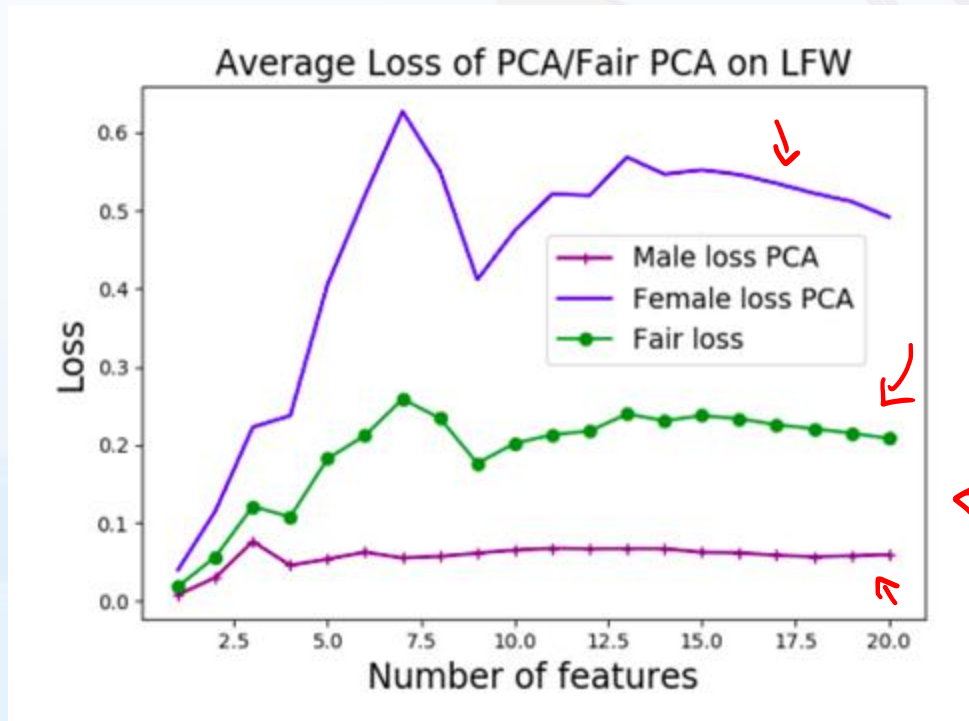
<https://sites.google.com/site/ssamadi/home/fair-pca-homepage>

Theorem *Let U be a solution to the Fair PCA problem (1), then*

$$\frac{1}{|A|} \text{loss}(A, \underline{U_A}) = \frac{1}{|B|} \text{loss}(\underline{B}, \underline{U_B}).$$

$$U = \begin{bmatrix} U_A \\ U_B \end{bmatrix}$$

Theorem 5.1. *There is a polynomial-time algorithm that outputs an approximation matrix of the data such that it is either of rank d and is an optimal solution to the fair PCA problem OR it is of rank $d + 1$, has equal losses for the two populations and achieves the optimal fair PCA objective value for dimension d .*





Fair Supervised learning

No single answer

Statistical parity
Group fairness
Demographic parity
Conditional statistical parity
Equal opportunity
Equalized odds
Conditional procedure accuracy equality
Disparate mistreatment
Balance for positive class
Balance for negative class
Predictive equality
Conditional use accuracy equality
Predictive parity
Calibration

Statistical Parity

X – set of all data points (people)

C – Set of all data points (people) belonging to a certain protected group

$M: X \rightarrow \{0,1\}$ - a classifier, say *logistic regression*

Bias of classifier for the protected group:

$$\text{Parity}(M,C) = \Pr(M(x) = 1/x \text{ in } C) - \Pr(M(x) = 1)$$

Ideal classifier \Rightarrow Parity = 0

In practice, $\min |\text{parity}(M,C)|$

Equality of Opportunity

X – set of all data points (people) ✓

C – Set of all data points (people) belonging to a certain protected group ✓

$M: X \rightarrow \{0,1\}$ - a classifier, say *logistic regression*. ✓

$$\text{opportunity_inequality}(M,C) = \Pr(M(x) = 1 \mid y=1 \text{ and } C) - \Pr(M(x) = 1 \mid y=1)$$

Ideal classifier: opportunity_inequality = 0

In practice: $\min \text{opp_ineq}(M,C)$

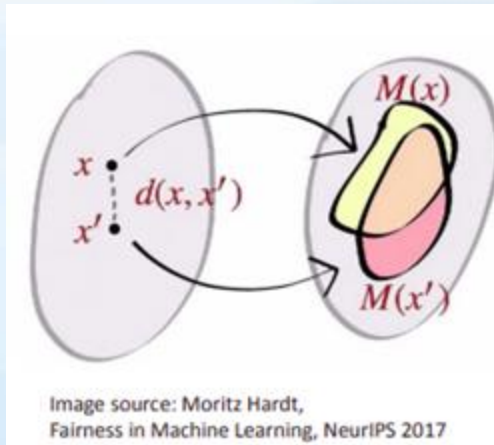
Individual Fairness

\mathbf{X} – set of all data points (people)

$M: \mathbf{X} \rightarrow \Delta$ (simplex) - a classifier, say *logistic regression*.

$d: \mathbf{X} * \mathbf{X} \rightarrow \mathbb{R}$ - distance function

$D: \Delta * \Delta \rightarrow \mathbb{R}$ – distance between probability vectors



$$D(M(x), M(x')) \leq d(x, x') \text{ for any } x, x'$$

Outcome Test (Predictive Parity)

X – set of all data points (people)

C – Set of all data points (people) belonging to a certain protected group

$M: X \rightarrow \{0,1\}$ - a classifier, say *logistic regression*.

$$\text{Pred_parity}(M,C) = \Pr(y=1 \mid M(x) = 1 \text{ and } C) - \Pr(y=1 \mid M(x) = 1)$$

Ideal classifier: $\text{Pred_parity} = 0$

In practice: $\min |\text{Pred_parity}(M,C)|$

Impossibility Result

$\text{Opportunity_inequality}(M,C) = \Pr(M(x) = 1 \mid y=1 \text{ and } C) - \Pr(M(x) = 1 \mid y=1)$

$\text{Opportunity_inequality}(M,C) = \Pr(M(x) = 1 \mid y=1 \text{ and } C) - \Pr(M(x) = 1 \mid y=1)$

Calibration within groups – $\Pr(M(x) = 1 \mid x \text{ in } C) \approx \Pr(y=1 \mid x \text{ in } C)$

Inherent Trade-Offs in the Fair Determination of Risk Scores

Jon Kleinberg *

Sendhil Mullainathan †

Manish Raghavan ‡

Fair Logistic Regression

<https://dl.acm.org/doi/fullHtml/10.1145/3308560.3317584>

$$\text{Parity} = |\Pr(M(x) = 1 | C = 1) - \Pr(M(x) = 1 | C = 0)|$$

Regularized logistic loss on dataset

minimize $f_D(\mathbf{w})$
 subject to $g_D(\mathbf{w}) \leq \tau, g_D(\mathbf{w}) \geq -\tau,$

equivalently.

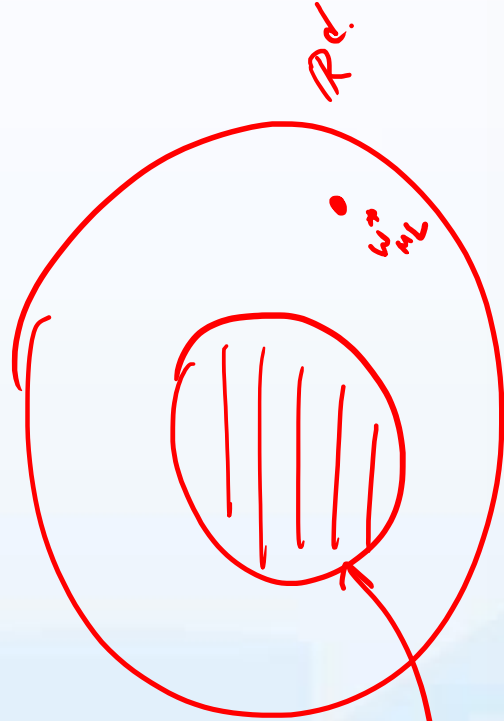
$$|g_D(\mathbf{w})| \leq \tau$$

$$g_D(\mathbf{w}) = \sum_{i=1}^n (s_i - \bar{s}) \mathbf{x}_i^T \mathbf{w}$$

Measures how correlated are the prediction probabilities to the Protected attribute.

Protected attribute for data point i

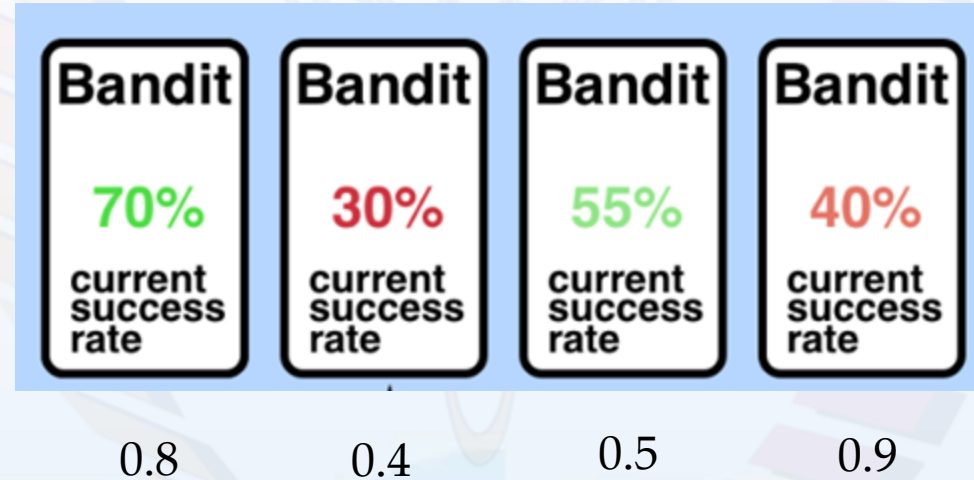
$$S = \begin{bmatrix} s_1 - \bar{s} \\ s_2 - \bar{s} \\ \vdots \\ s_n - \bar{s} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad \tau$$



$$\left\{ \frac{\sum_{i=1}^n (s_i - \bar{s}) \mathbf{x}_i^T \mathbf{w}}{\sqrt{\sum_{i=1}^n (s_i - \bar{s})^2}} \right\} \leq \tau$$

τ

Fair Multi Armed Bandits



Case in point: Swiggy/Zomato wants to assign partners to orders.
Goal: Maximize expected reward over time.

A simple strategy: Pick current best arm with (0.9) probability and uniformly at random with 0.1 prob.

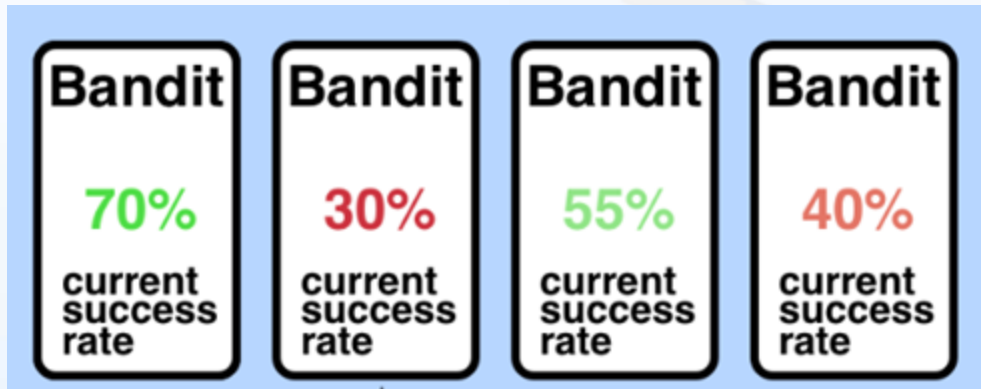
Fair Multi Armed Bandits

Notion of Fairness:

If partner-a is **truly** better than partner b,
then

Probability that algorithm assigns partners a \geq
Probability that algorithm assigns partners b

~~A simple strategy: Pick current best arm with (0.9) probability and uniformly at random with 0.1 prob.~~

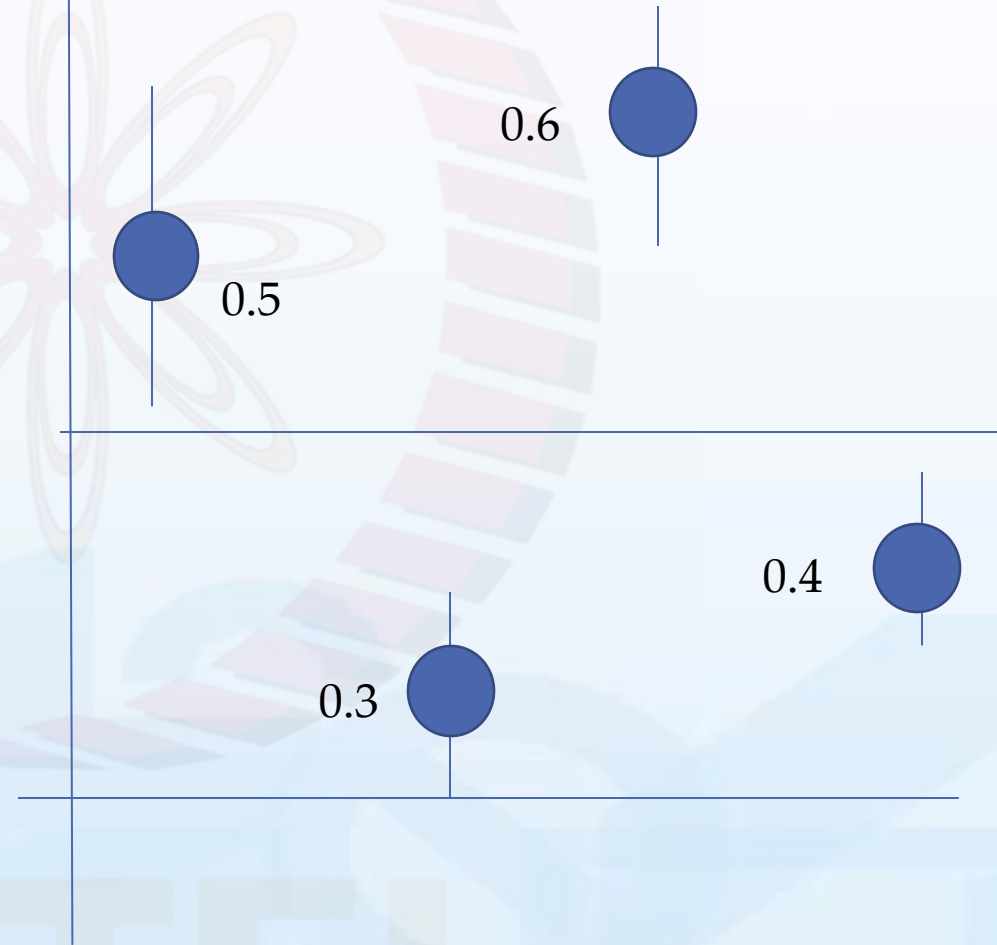


0.5


0.35

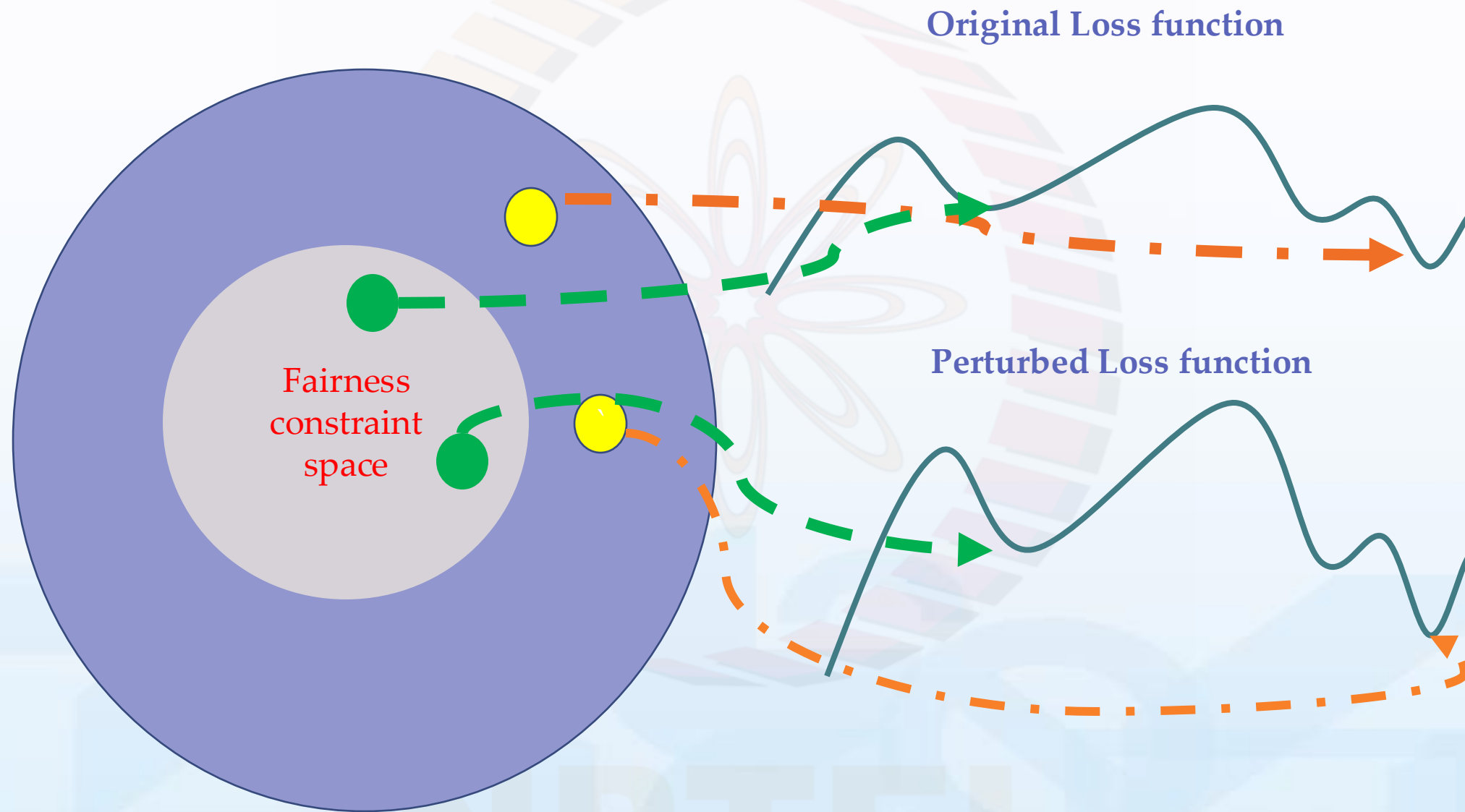
0.6

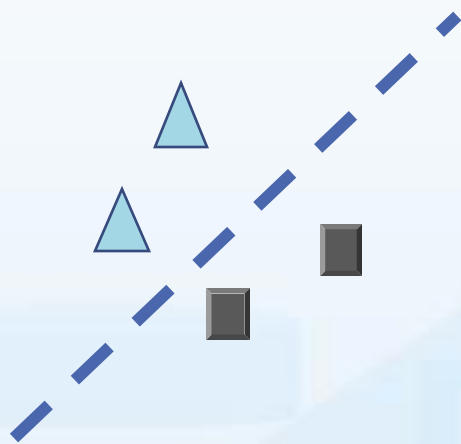
0.38



IDEA: Sample uniformly from those whose confidence interval overlaps with the winner

- 
- Fair PCA
 - Fair LogReg
 - Fair MAB





Model

“

If you can't explain
something to a first year
student, then you haven't
really understood.

—
RICHARD FEYNMAN

NPTTEL

Thank you!



<https://cerai.iitm.ac.in>

arunr@cse.iitm.ac.in

NPTTEL